

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 11-007447

(43)Date of publication of application : 12.01.1999

(51)Int.Cl. G06F 17/27
G10L 3/00

(21)Application number : 09-160954

(71)Applicant : NIPPON TELEGR & TELEPH CORP <NTT>

(22)Date of filing : 18.06.1997

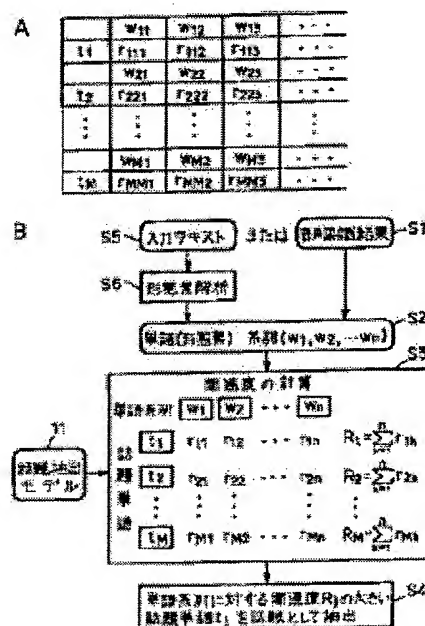
(72)Inventor : OFU KATSUTOSHI
MATSUOKA TATSUO
MATSUNAGA SHOICHI

(54) TOPIC EXTRACTING METHOD, TOPIC EXTRACTION MODEL TO BE USED FOR THE EXTRACTING METHOD, PREPARING METHOD FOR THE TOPIC EXTRACTION MODEL, AND TOPIC EXTRACTION PROGRAM RECORDING MEDIUM

(57)Abstract:

PROBLEM TO BE SOLVED: To appropriately extract a topic (heading) expressing the contents of a continuous voice.

SOLUTION: A continuous large vocabulary voice is recognized (S1) and a word sea prepared (S2) by using a model 11 stored by executing syntax analysis for the headings and texts of many newspaper topics, obtaining respective topic words of the headings and words in the texts, finding out the appearance frequency of respective words and the cooccurrence frequency of combination of each topic word and a text word in the same news item, and finding out the degree of association between the topic word and the text word by mutual information volume or an x2 method. The degree of association between each topic word and each word in the word sequence is found out by the model 11 to prepare a association degree sequence, the sum of association degrees in respective association degree sequences is found out (S3) and a topic word corresponding to the maximum sum is outputted (S4).



(51) Int.Cl.⁶

G 0 6 F 17/27

G 1 0 L 3/00

識別記号

5 6 1

F I

G 0 6 F 15/38

G 1 0 L 3/00

N

5 6 1 G

審査請求 未請求 請求項の数12 O L (全 6 頁)

(21) 出願番号 特願平9-160954

(22) 出願日 平成9年(1997) 6月18日

(71) 出願人 000004226

日本電信電話株式会社

東京都新宿区西新宿三丁目19番2号

(72) 発明者 大附 克年

東京都新宿区西新宿三丁目19番2号 日本
電信電話株式会社内

(72) 発明者 松岡 達雄

東京都新宿区西新宿三丁目19番2号 日本
電信電話株式会社内

(72) 発明者 松永 昭一

東京都新宿区西新宿三丁目19番2号 日本
電信電話株式会社内

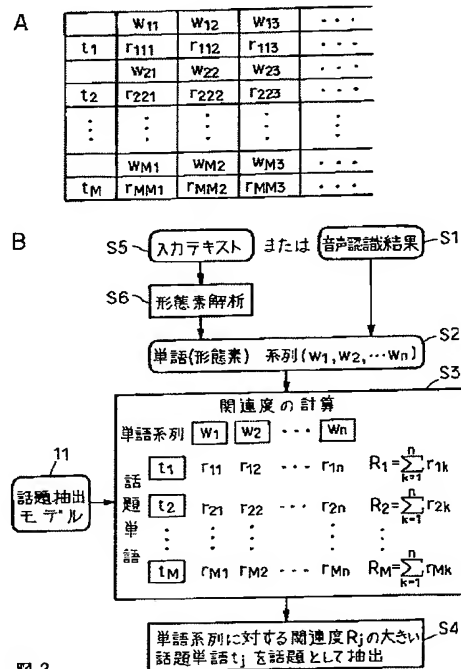
(74) 代理人 弁理士 草野 卓

(54) 【発明の名称】 話題抽出方法及びこれに用いる話題抽出モデルとその作成方法、話題抽出プログラム記録媒体

(57) 【要約】

【課題】 連続音声の内容を表わす話題（見出し）を適切に抽出する。

【解決手段】 大量の新聞記事の見出しと本文とを形態素解析し、その見出しの各話題単語と文中単語を得、その各出現頻度と、同一記事で話題単語と文中単語の組合せの共起頻度を求め、相互情報量又は χ^2 法により話題単語と文中単語との関連度を求めて格納したモデル11を用い、連続大語彙音声を生声認識し（S1）、単語系列を作り（S2）、各話題単語と単語系列の各単語との関連度をモデル11が求めて、関連度系列を作り各関連度系列における関連度の和を求め（S3）、その和の最大のもので対応する話題単語を出力する（S4）。



【特許請求の範囲】

【請求項 1】 複数の単語の系列の内容を表す話題単語を抽出するために用いられるモデルであって、複数の話題単語と、その各話題単語と、各単語との関連度とがそれぞれ格納されている話題抽出モデル。

【請求項 2】 話題単語と各単語との関連度は、話題単語と複数単語との関連度であることを特徴とする請求項 1 記載の話題抽出モデル。

【請求項 3】 上記関連度は話題単語と、各単語との相互情報量に基づくものであることを特徴とする請求項 1 又は 2 記載の話題抽出モデル。

【請求項 4】 上記関連度は話題単語と各単語との x^2 ベクトル法にもとづくものであることを特徴とする請求項 1 記載の話題抽出モデル。

【請求項 5】 請求項 1 乃至 4 の何れかに記載した話題抽出モデルを用いて入力された複数の単語の系列の内容を表す話題単語を抽出する方法であって、上記話題抽出モデル中の各話題単語ごとに、これと上記入力単語系列の各単語との関連度を上記話題抽出モデルを参照して求めて関連度系列をそれぞれ作り、これら各関連度系列の各関連度の和を求めて上記単語系列に対する各話題単語の関連度を求め、これら単語系列に対する関連度中の大きいものから順に Q 個 (Q は 1 以上の整数) のものとそれぞれ対応する話題単語を出力することを特徴とする話題抽出方法。

【請求項 6】 上記関連度系列の各関連度に対し、これと対応する単語の尤度で重み付けて上記各関連度の和を求め、その結果を特徴とする請求項 5 記載の話題抽出方法。

【請求項 7】 連続した音声信号を単語音声認識して、上記入力単語系列を得ることを特徴とする請求項 5 又は 6 記載の話題抽出方法。

【請求項 8】 上記認識結果として複数の上位の候補系列を上記入力単語系列とすることを特徴とする請求項 7 記載の話題抽出方法。

【請求項 9】 入力テキストを、形態素解析し、その解析結果の形態素を上記入力単語系列とすることを特徴とする請求項 5 又は 6 に記載の話題抽出方法。

【請求項 10】 本文とその見出しよりなる多数のテキストを学習データとし、この学習データの見出し、本文をそれぞれ形態素解析して、見出しの形態素としての話題単語と、本文の形態素としての文中単語を得る工程と、

上記各話題単語の出現頻度と、上記各文中単語の出現頻度と、1 つのテキスト中の上記話題単語と上記文中単語の各組み合わせが同時に得られる共起頻度とをそれぞれ計数する工程と、

上記話題単語の出現頻度と文中単語の出現頻度と各共起頻度とを用いて各話題単語と各文中単語との関連度を求めて話題抽出モデルを得る工程とを有する話題抽出モデル作成方法。

【請求項 11】 上記出現単語中の出現頻度が所定値以下のものを省略し、上記文中単語中の出現頻度の順位が所定値以下のものを省略し、上記出現単語及び上記文中単語中の情報検索という観点から意味的情報を伝達する名詞・動詞などの品詞のもの以外を省略し、かつ 1 つのテキスト中に出現する上記話題単語及び文中単語の組み合わせが所定回数以下の組み合わせを省略して残りの話題単語及び文中単語を用いて上記関連度を求めることを特徴とする請求項 10 記載の話題抽出モデル作成方法。

【請求項 12】 入力音声を連続音声認識して入力単語系列を得、複数の話題単語と、その各話題単語と、各単語との関連度とがそれぞれ格納された話題抽出モデルを参照して、上記話題単語ごとにこれと上記入力単語系列中の各単語との関連度を求めて関連度系列を得、上記各関連度系列から、上記各話題単語ごとの上記入力単語系列の関連度を求め、これら入力単語系列の関連度中の関連度が最大のものから順に Q 個 (Q は 1 以上の整数) のものとそれぞれ対応する話題単語を上記入力単語系列の内容を表わす話題として出力することをコンピュータを用いて行うためのプログラムを記録した記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】この発明は、連続発声された音声の単語認識結果やテキストを形態素解析により分割された単語などの単語列に対し、その内容を表わす話題を抽出する方法、その話題抽出に用いる話題抽出モデルとそのモデルの作成方法に関する。

【0002】

【従来の技術】連続発声された音声からのその内容を表わす話題抽出では 5 ～ 10 種類の分野のうちのいずれかの分野に依存度の高いキーワードを予め選択しておき、それらのキーワードを音声区間中から検出 (キーワードスポッティング) して、検出されたキーワード集合が最も高い依存度を示す話題を結果として出力する方法が多くとられる。例えば横井、河原、堂下、"キーワードスポッティングに基づくニュース音声の話題同定"、情処研報、SLP6-3、1995。櫻井、有木、"キーワードスポッティングによるニュース音声の索引付けと分類"、信学技法、SP96-66、1996。R.C.Rose, E.L.Chang, and R.P.Lippmann, "Techniques for Information Retrieval from Voice Messages", Proc.ICASSP-91, pp.317-320, 1991。などに示されている。

【0003】また従来の文章 (テキスト) から話題を抽出する方法は文中の特定の個所を抽出して行うもので、その処理が複雑であった。

【0004】

【発明が解決しようとする課題】従来の連続音声の話題抽出方法では、限られた数のキーワードしか用いること

10

20

30

40

50

ができず、またキーワードの数を増やした場合には、誤って検出されるキーワードが増えてしまう、また話題の分野が少ないため、情報検索や索引付けに使うことができないという問題があった。また従来のテキストよりの話題抽出方法は、特定の個所を探して行うため、その処理が複雑であった。これを連続音声の話題抽出に適用すると、その所定個所についての単語認識を誤ると、話題抽出は誤ってしまう。

【0005】この発明の目的は比較的簡単な処理で話題を抽出することができる話題抽出方法、そのプログラムを記録した記録媒体と、上記話題抽出に用いる話題抽出モデルと、その作成方法を提供することにある。

【0006】

【課題を解決するための手段】この発明の話題抽出モデルは、本文とその見出しよりなるテキストを多数用いて、それぞれ形態素解析を用い、本文単語と話題単語（見出し中の）と得、これら本文単語の出現頻度、話題単語の出現頻度、1テキスト中に本文単語と話題単語の組み合わせが同時に存在する共起頻度をそれぞれ求め、これら頻度から各話題単語と、本文単語との関連度を求め、これらを話題抽出モデルとして格納しておく。

【0007】この発明の話題抽出方法では前記この発明の話題抽出モデルを用い、入力音声の音声認識や入力テキストの形態素解析で、入力単語系列を得、各話題単語と入力単語系列中の各単語との関連度とを話題抽出モデルを参照して求めて話題単語ごとの関連度系列を得、これら関連度系列から各話題単語の入力単語系列との関連度をそれぞれ求め、これら入力単語系列の関連度中の大きいものと対応する話題単語を入力音声又はテキストに対する話題として出力する。

【0008】この発明の記録媒体はこの発明の話題抽出方法をコンピュータで実行させるためのプログラムが記録されている。

【0009】

【発明の実施の形態】まずこの発明の話題抽出モデルと*

$$I(w_i : t_j) = \log(P(w_i, t_j) / P(w_i) P(t_j)) \dots (1)$$

$P(w_i, t_j)$: w_i と t_j が同時に出現する確率
 $P(w_i)$: w_i の出現確率、 $P(t_j)$: t_j の出現確率

χ^2 法に基づく関連度

$$\chi_{ij}^2 = (f_{ij} - F_{ij})^2 / F_{ij}$$

【0012】

【数1】

$$F_{ij} = \frac{\sum_{\ell=1}^M f_{i\ell}}{\sum_{k=1}^N \sum_{\ell=1}^M f_{k\ell}} \cdot \sum_{k=1}^N f_{kj}$$

※

$$I'(w_i : t_j) = I(w_i : t_j), P(w_i, t_j) \neq 0 \text{ の場合}$$

* この作成方法の実施例を説明する。話題抽出モデルの学習（作成）はある話題について述べられているテキストとその内容を表わす複数の話題単語との組を大量に用いて行う。一例として新聞記事の本文と見出しを用いて話題抽出モデルを学習（作成）する場合、約5年分の新聞記事よりその見出しと本文とをそれぞれ取出し（S1）、これらを形態素解析を行い（S2）、単語（形態素）に分割し、見出しの形態素（話題単語）と、本文の形態素（文中単語）とを得る。

10 【0010】これら話題単語と文中単語について、大量のデータにおける出現頻度と、共起頻度とを用いて、文中単語と話題単語との関連度を求める。しかし、文中単語と話題単語の組み合わせは非常に莫大な数になる。従ってこの実施例では話題単語については、出現回数が2回以上の単語に限り（S3）、文中単語については出現頻度が上位15万の単語のみを選出し（S4）、更に情報検索という観点からより意味的情報を伝達すると考えられる名詞、動詞などの内容語に着目し、ここでは話題単語、文中単語の何れについても名詞、動詞、形容詞、

20 形容動詞、副詞のみを取出す（S5）。更に話題単語と文中単語との組合せで同一記事に出現するのが1度しかなかったものは除外し、つまり話題単語と文中単語の組み合わせで同一記事に出現することが2回以上のもののみを残した（S6）。このようにして話題単語の総頻度 12.3×10^6 が 6.3×10^6 となり総数 136×10^3 が 74×10^3 となり、文中単語の総頻度 218.8×10^6 が 90.1×10^6 となり総数 640×10^3 が 147×10^3 となり、2回以上起きた共起の組み合わせは約5800万種類となった。

30 【0011】この約5800万種類について、これら単語の出現頻度と共起頻度を用いて文中単語と話題単語との関連度を求める。文中単語 w_i と話題単語 t_j との関連度は以下のようにして求める。相互情報量に基づく関連度

※ N : 文中単語の種類数、M : 話題単語の種類数、
 f_{ij} : 話題単語 t_j に対する文中単語 w_i の頻度、
 F_{ij} : 話題単語 t_j に対する文中単語 w_i の理論（期待）度数

40 相互情報量の計量において、学習データ中に文中単語 w_i と話題単語 t_j の共起が観測されない場合、 $P(w_i, t_j) = 0$ となり、関連度の合計を求める際に問題が生じる。そこで、共起が観測されなかった場合には情報が得られなかったものとして、実際には次式のように相互情報量に基づく関連度を計算する。

【0013】

一方、 x^2 法における理論度数 F_{ij} とは、全ての話題単語に対して文中単語 w_i が等確率で出現した場合の出現頻度である。実際の出現頻度と理論度数とのずれが大きければ、その文中単語はその話題単語に対して偏って出現していることになる。しかし、上述の x^2 法の式では、実際の出現頻度 f_{ij} が理論度数 F_{ij} より小さい場合にも、関連度が正の値となってしまうため、実際には次式のように x^2 法に基づく関連度を計算する。

【0014】

$$x_{ij}^2 = x_{ij}^2, \quad f_{ij} - F_{ij} \geq 0 \text{ の場合} \\ 0, \quad f_{ij} - F_{ij} < 0 \text{ の場合}$$

従って、ステップ S 6 で得られた文中単語 w_i と話題単語 t_j との各組み合わせについて、その各頻度 P

$(w_i) : P(t_j), P(w_i, t_j)$ 、または f_{ij} をそれぞれ演算し (S 7)、頻度テーブル 11 に格納する。これを学習データが終るまで行う (S 8)。学習データが終ると、頻度テーブル 11 内に演算した頻度を用いて関連度 $I(w_i, t_j)$ 又は F_{ij} の計算を行って話題抽出モデルを得る (S 9)。

【0015】従って話題抽出モデルは図 2 A に示すように、話題単語の種類 t_1, t_2, \dots, t_k それぞれについて、これと 2 回以上共起する文中単語、つまり t_1 については $w_{11}, w_{12}, w_{13}, \dots$ との関連度 $r_{11}, r_{12}, r_{13}, \dots$ が、また t_2 については $w_{21}, w_{22}, w_{23}, \dots$ との関連度 $r_{21}, r_{22}, r_{23}, \dots$ が、以下同様に文中単語との関連度が格納されている。

【0016】次にこの話題抽出モデルを用いて連続入力単語列から話題を抽出する方法を図 2 B を参照して説明する。連続発声される音声を入力とする場合、その入力音声を単語音声認識 (S 1)、認識結果として単語系列 w_1, w_2, \dots, w_n を得る (S 2)、これら単語系列 w_1, w_2, \dots, w_n の各単語について、話題抽出モデル*

$$I(x_1 : x_2 : \dots : x_n) = \log \frac{[\prod P(x_i, x_j)] \cdot [\prod P(x_i, x_j, x_k, x_l)]}{[\prod P(x_i)] \cdot [\prod P(x_i, x_j, x_k)]}$$

$P(w_i, t_j) = 0$ の場合

* 11 を参照して、その各話題単語 t_1, t_2, \dots, t_k に対する関連度を求める。つまり認識単語 w_i に対する話題単語 t_1, t_2, \dots, t_k との各関連度 $r_{11}, r_{21}, \dots, r_{k1}$ を求め、単語 w_2 に対する話題単語 t_1, t_2, \dots, t_k との各関連度 $r_{12}, r_{22}, \dots, r_{k2}$ を求め、以下同様に求める。

【0017】各話題単語 t_1, t_2, \dots, t_k についての各認識単語 w_1, w_2, \dots, w_n との関連度の合計、

つまり単語系列に対する関連度 R_j を計算する。即ち、話題単語 t_1 については $r_{11}, r_{12}, \dots, r_{1n}$ の和 $R_1 = \sum_{k=1}^n r_{1k}$ を求め、 t_2 については $r_{21}, r_{22}, \dots, r_{2n}$ の和 $R_2 = \sum_{k=1}^n r_{2k}$ を求め、以下同様に R_3, \dots, R_k を求める (S 3)。これら単語系列に対する関連度 R_1, \dots, R_k 中で関連度が大きいものから順に Q 個 (Q は 1 以上の整数) のものとそれぞれ対応する話題単語 t_j の具合を、その単語系列に対する話題とする (S 4)。 Q は 1 でもよいが、通常は複数で例えば 5 程度である。関連度 R_1, \dots, R_k 中の大きいものから順に対応する話題単語の複数列を候補として出力してもよい。

【0018】単語系列から話題の抽出としてはテキストを入力し (S 5)、これを形態素解析し (S 6)、形態素つまり単語列 w_1, w_2, \dots, w_n を得て、これを音声入力の場合と同様に話題抽出モデル 11 を用いて処理して、テキストに対する話題を抽出することもできる。関連度を w_i と t_j の相互情報量に基づいて求める場合は式 (1)、つまり 2 点間の相互情報量に基づいて決めた。一方、 n 点間の相互情報量は次式で定義される。

【0019】

【数 2】

... (2)

Π は、あい異なる添字の全ての組み合わせについて計算する。従って x_1, x_2, \dots, x_n 中 1 つの話題単語と他の $n-1$ 個を文中単語との相互情報量を $I(x_1 : x_2 : \dots : x_n)$ により求めることができる。このように、複数の文中単語と 1 つの話題単語との関連度を求めておくと、例えば「コンピュータ」と「インターネット」の関連度、また「ネットワーク」と「インターネット」の関連度はそれ程大きくないが、「コンピュータ」と「ネットワーク」が同じ文中にあった場合の「インターネット」への関連度が大きくなるような話題抽出モデルの学習ができる。つまり式 (1) の関連度では話題と

して「インターネット」を抽出できない場合に、式 (2) の関連度によると「インターネット」を話題として抽出でき、適切な話題を抽出することができるがある。

【0020】話題単語 t_k と単語系列 w_1, w_2, \dots, w_n との関連度 R_k は t_k に対する各単語の関連度の和 $r_{k1} + r_{k2} + \dots + r_{kn}$ で求められるが、その加算の際に各単語に対する重み s_1, s_2, \dots, s_n をそれぞれ付けて、 $r_{k1} \times s_1 + r_{k2} \times s_2 + \dots + r_{kn} \times s_n$ というようにして、より適切な関連度 R_k を得るようにすることもできる。ここで重み s_1, s_2, \dots, s_n

としては、各単語 w_1, w_2, \dots, w_n のその音声認識時の単語の確からしさ（音響的尤度）や言語的尤度、つまりその単語がその前の単語に対し、文法や言語上存在する確からしさ（大語彙音声認識に用いられる言語モデルに示されている）を用いることができる。

【0021】音声認識結果の単語系列に対して話題抽出を行う際に、認識言語系列候補の第1位だけでなく、上位 b 位までの候補（ $w_{1-1}, w_{1-2}, \dots, w_{1-n1}$ ）,（ $w_{2-1}, w_{2-2}, \dots, w_{2-n2}$ ） \dots （ $w_{b-1}, w_{b-2}, \dots, w_{b-nb}$ ）を用いて話題抽出を行う。この際、順位の高い程重みが大きくなるようにすることもできる。この場合第1位から第 b 位までの候補系列は、相互に1単語又は2単語しか違いがない候補系列が多くなる。よってこれら候補系列を、その同一単語を排除して複数単語木構造乃至単語ネットワークあるいは単語ラティスの配列とし、これを用いて第1位～第 b 位の候補系列について話題抽出をするようにすれば、その複数の候補系列を少ない容量のメモリに格納して処理することができる。

【0022】

【発明の効果】評価は、ニュース音声の書き起こし文および2万語彙の大語彙連続音声認識システムによる音声認識結果に対してこの発明の評価を行った。書き起こし文に対して3人の被験者が人手で付与した話題を評価対象とした。話題抽出モデルが出力した話題単語のうち上位5単語までを出力結果とした場合の適合率（抽出した話題単語のうち、正解の話題単語の割合）は、3人の被験者の付与した話題に対して70%以上となった。また、単語誤り率25%の音声認識結果に対する話題抽出結果の適合率も65%以上となった。各被験者が付与した話題間の重複は約70%であるので、この話題抽出結果は利用可能な精度であるといえる。関連度の尤度とし

て χ^2 法を用いた方が相互情報量を用いた場合より良い結果が得られた。

【0023】この発明によれば、大量のテキストデータを用いて非常に多くの文中単語および話題単語間の関連度を学習した話題抽出モデルを用いることにより、テキストおよび誤りを含む大語彙連続音声認識結果から詳細な話題抽出を行うことができるという利点がある。つまり、音声からの話題抽出において、連続音声認識技術を用いることにより、限られた数のキーワードを検出するキーワードスポッティングに基づく方法に比べ、音声中の多くの情報を用いて話題抽出を行うことが可能であり、また、音声の内容を表す単語（話題単語）を複数抽出することにより、音声をいくつかの分野に分類する話題抽出（話題同定・話題認識）に比べ、詳細な話題が抽出できるという利点がある。

【0024】特に従来のテキストに対する話題抽出では、特定の関係のものを抽出するため、複雑な処理を必要としたが、この発明では比較的簡単に行うことができる。特に連続音声に対する抽出ではその特定部分に対して認識誤りが生じると、致命的であるが、この発明は文全体の各単語に対して関連性をみるため正しく話題を抽出することができる。

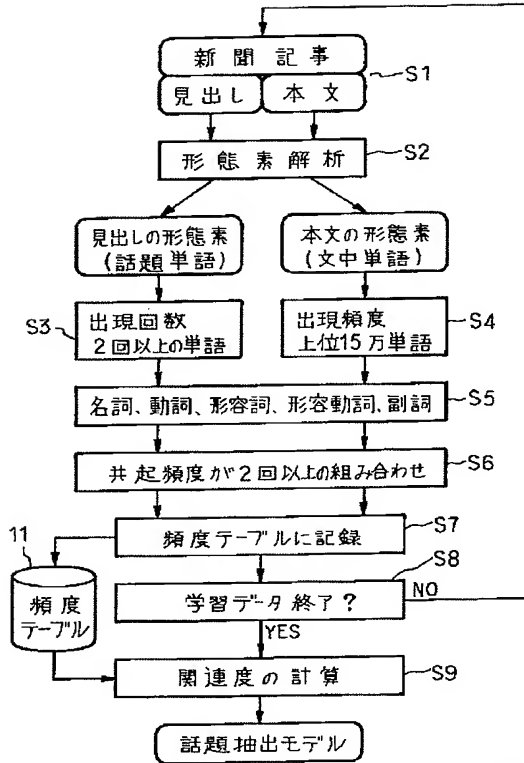
【0025】またこのような正しい抽出ができるのは、大量のテストデータを用いて作成した話題単語と各単語との関連度を記憶した話題抽出モデルを用いるからである。

【図面の簡単な説明】

【図1】この発明のモデル作成方法を示す流れ図。

【図2】Aはこの発明の話題抽出モデルの例を示す図、Bはこの発明の話題抽出方法を示す図である。

【図 1】



【図 2】

